



Empirical generalization assessment of neural network models

Larsen, Jan; Hansen, Lars Kai

Published in:

Proceedings of the 1995 IEEE Workshop on Neural Networks for Signal Processing

Link to article, DOI:

[10.1109/NNSP.1995.514876](https://doi.org/10.1109/NNSP.1995.514876)

Publication date:

1995

Document Version

Publisher's PDF, also known as Version of record

[Link back to DTU Orbit](#)

Citation (APA):

Larsen, J., & Hansen, L. K. (1995). Empirical generalization assessment of neural network models. In *Proceedings of the 1995 IEEE Workshop on Neural Networks for Signal Processing* (pp. 30-39). IEEE. <https://doi.org/10.1109/NNSP.1995.514876>

General rights

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the public portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the public portal

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

EMPIRICAL GENERALIZATION ASSESSMENT OF NEURAL NETWORK MODELS

Jan Larsen and Lars Kai Hansen
The Computational Neural Network Center (CONNECT)
Electronics Institute, Building 349
Technical University of Denmark
DK-2800 Lyngby, Denmark
emails: jlarsen,lkhansen@ei.dtu.dk

Abstract. This paper addresses the assessment of generalization performance of neural network models by use of empirical techniques.

We suggest to use the cross-validation scheme combined with a resampling technique to obtain an estimate of the generalization performance distribution of a specific model. This enables the formulation of a bulk of new generalization performance measures. Numerical results demonstrate the viability of the approach compared to the standard technique of using algebraic estimates like the FPE [1].

Moreover, we consider the problem of comparing the generalization performance of different competing models. Since all models are trained on the same data, a key issue is to take this dependency into account.

The optimal split of the data set of size N into a cross-validation set of size $N\gamma$ and a training set of size $N(1-\gamma)$ is discussed. Asymptotically (large data sets), $\gamma_{\text{opt}} \rightarrow 1$ such that a relatively larger amount is left for validation.

INTRODUCTION

Consider a tapped-delay neural network predicting a stochastic output signal¹ $y(k)$ from observations of the L -dimensional stochastic input vector signal $x(k) = [x(k), x(k-1), \dots, x(k-L+1)]^\top$. Let $f(\cdot)$ denote the mapping of the neural network, and w the vector of network weights (parameters). Then the prediction is given as: $\hat{y}(k) = f(x(k); w)$.

¹For convenience, we focus on single output models.

The basic object of interest in neural network modeling is the conditional input-output distribution $p(y|\mathbf{x})$, i.e., the probability distribution of the output conditioned on a test input vector, see e.g., [16]. Normally the network is trained to implement the conditional mean², $E\{y|\mathbf{x}\} = \int y \cdot p(y|\mathbf{x}) dy$. The first source of uncertainty is the inherent prediction error $\varepsilon = y - E\{y|\mathbf{x}\}$ which – per definition – cannot be modeled. Another considerable source of uncertainty is the estimation of $E\{y|\mathbf{x}\}$ from a limited number of training data.

This paper deals with empirical assessment of model quality expressed in terms of generalization performance defined as prediction accuracy on future data. Reliable estimates of the generalization performance of a particular model is very important for practical applications. Moreover, in order to choose the best model from a pool of candidate model architectures³, one requires a test which determines if a particular model has a significantly higher generalization performance than a competing model. The empirical framework enables both *absolute* and *comparative* generalization assessment.

The generalization performance can be decomposed into three components, see e.g., [3], [6]. The first term is due to the inherent prediction error, ε . The second term expresses the insufficiency of the neural architecture⁴ to model the conditional mean, and is often referred to as the model bias. Finally, the third term reflects finite training set effects, also known as the model variance. While the first term – per definition – cannot be decreased, there will normally exist a trade off between bias and variance which is accomplished by optimizing the architecture, e.g., by using pruning techniques.

ON GENERALIZATION PERFORMANCE

Suppose the network is trained by minimizing a cost function, viz. the sum of a loss function, $S_N(\mathbf{w})$, and a regularization term $R(\mathbf{w})$, i.e.,

$$C_N(\mathbf{w}) = S_N(\mathbf{w}) + R(\mathbf{w}) = \frac{1}{N} \sum_{k=1}^N \ell(y(k), \hat{y}(k); \mathbf{w}) + R(\mathbf{w}) \quad (1)$$

where $\ell(\cdot)$ measures the distance between the output $y(k)$ and the network prediction $\hat{y}(k) = f(\mathbf{x}(k); \mathbf{w})$. Even though much of the material in this paper applies for general loss functions, often the mean square error loss function, $\ell = (y - \hat{y})^2 = e^2$ is considered. N defines the number of training examples, i.e., input-output pairs of the training set: $\mathcal{D} = \{(\mathbf{x}(k), y(k))\}_{k=1}^N$.

Training on the full set of examples provides the estimated weight vector $\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} C_N(\mathbf{w})$. The generalization error, G , is defined as the *expected*

²This is optimal when using a mean square error cost function, see e.g., [16].

³E.g., feed-forward neural nets with different input lag-space and number of hidden units.

⁴The architecture is presumed to be finite, i.e., the weight vector is finite dimensional.

loss of the estimated model on a test sample (x, y) independent of those in the training set,

$$G(\hat{\mathbf{w}}) = E\{\ell(y, \hat{y}; \hat{\mathbf{w}})\} = \int \ell(y, \hat{y}; \hat{\mathbf{w}}) \cdot p(x, y) dx dy \quad (2)$$

where $E\{\cdot\}$ denotes expectation w.r.t. the unknown joint input-output probability density $p(x, y)$. $G(\hat{\mathbf{w}})$ depends on the actual training set \mathcal{D} through the estimated weights $\hat{\mathbf{w}}$ and has the lower bound $G_{\min} = G(\mathbf{w}^*)$. \mathbf{w}^* denotes the optimal weight vector $\mathbf{w}^* = \arg \min_{\mathbf{w}} E\{C_N(\mathbf{w})\} = \arg \min_{\mathbf{w}} [G(\mathbf{w}) + R(\mathbf{w})]$ which corresponds to training on an infinite training set. Under fairly mild assumptions, it is possible to show $\lim_{N \rightarrow \infty} \hat{\mathbf{w}} = \mathbf{w}^*$, see e.g., [10], [11], [17]. G_{\min} expresses the fundamental uncertainty of y when x is known, and furthermore the potential lack of modeling capability, i.e., the network is incapable of implementing the optimal⁵ function, $g(x) \doteq \arg \min_{\phi(x)} E\{\ell(y, \phi(x))\}$, $\phi(x) : \mathbb{R}^p \mapsto \mathbb{R}$. Insufficient modeling capability is due to two facts:

- In general, when using a finite architecture the model is incomplete, i.e., $f(x; \mathbf{w}^o) \neq g(x)$ where $\mathbf{w}^o = \arg \min_{\mathbf{w}} G(\mathbf{w})$ is the weights minimizing the expected loss using the architecture embodied by $f(\cdot)$.
- Regularization implies that the optimal weight vector \mathbf{w}^* does not equal \mathbf{w}^o ; even when using a complete model.

Since the N samples in \mathcal{D} are randomly selected from the joint density $p((x(1), y(1)), \dots, (x(N), y(N)))$ the generalization error $G(\hat{\mathbf{w}})$ is stochastic with a certain *generalization error probability distribution* $P(G) = \text{Prob}\{G(\hat{\mathbf{w}}) < G\}$ and associated density $p(G)$.

The object of interest for model design could be either the full generalization distribution or just the generalization error $G(\hat{\mathbf{w}})$ on the particular training set available. These cases are treated separately in the following. If one has a strong belief in the training set (e.g., if it is large) one might address $G(\hat{\mathbf{w}})$. Otherwise, it might be better to consider the training set as a typical set drawn randomly from the joint input-output distribution in order to reveal the generic characteristics of the employed model.

Since $p(G)$ depends on the true distribution of data, the model architecture, and the number of training data, it is impossible to fully characterize it. However, it is possible to give some general properties. Obviously, $p(G) = 0$ as $G < G_{\min}$. For finite training sets, $p(G)$ will have non-zero values for $G \geq G_{\min}$, and since $\lim_{N \rightarrow \infty} \hat{\mathbf{w}} = \mathbf{w}^*$, $p(G)$ tends to a Dirac delta function $\delta(G - G_{\min})$ for $N \rightarrow \infty$. If the model is complete, the loss function is the mean square error and no regularization is employed, it is possible to show⁶ asymptotically as $N \rightarrow \infty$, $G(\hat{\mathbf{w}}) \sim \sigma_\varepsilon^2(1 + \chi^2(p)/N)$ where σ_ε^2 is the prediction error noise variance and $\chi^2(p)$ is the χ^2 -distribution with $p = \dim(\mathbf{w})$ degrees of freedom.

⁵With respect to the employed loss function.

⁶This is done by using second order expansions of $G(\hat{\mathbf{w}})$ around \mathbf{w}^* , and the fact that the fluctuations $\Delta \mathbf{w} = \hat{\mathbf{w}} - \mathbf{w}^*$ are asymptotically Gaussian distributed. See e.g., [6], [7], [8] and [16].

The literature on generalization theory and estimation of generalization error does not in general address the problem of characterizing the full prediction risk probability density. Most work has focused on simple measures of location such as the *average generalization error*

$$\text{avr}(G) = E_{\mathcal{D}} \{G(\hat{\mathbf{w}})\} = \int G \cdot p(G) dG. \quad (3)$$

This includes algebraic estimators like FPE [1], FPER [7], GEN [5], GPE [9] and NIC [10] which are valid asymptotically $N \rightarrow \infty$ and make several assumptions on the model and statistics of the data. However, also algebraic estimates of fractiles of $p(G)$ have been developed, see e.g., [14], [15]. Thus the $1 - \alpha$ fractile $G_{1-\alpha}$ defined by $\text{Prob}\{G \leq G_{1-\alpha}\} = 1 - \alpha$ guarantees that the probability of G exceeding $G_{1-\alpha}$ is α , which can be set to some low percentage.

EMPIRICAL GENERALIZATION ERROR ESTIMATION

If the object of interest is the generalization error $G(\hat{\mathbf{w}})$ for the particular training set available, we consider the hold-out cross-validation technique [13] for estimating $G(\hat{\mathbf{w}})$. Suppose that a cross-validation set \mathcal{C} of $N_c = \lceil N\gamma \rceil$,⁷ $0 < \gamma < 1$, samples are hold out for cross-validation and denote by \mathcal{T} the remaining $N_t = N - N_c$ data for training, i.e., let $\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} C_{N_t}(\mathbf{w})$. The cross-validation estimate of $G(\hat{\mathbf{w}})$ then reads:

$$\hat{G}(\hat{\mathbf{w}}) = \frac{1}{N_c} \sum_{k \in \mathcal{C}} \ell(y(k), \hat{y}(k); \hat{\mathbf{w}}) \quad (4)$$

Under suitable regularity conditions, $\hat{G}(\hat{\mathbf{w}}) \rightarrow G(\hat{\mathbf{w}})$ as $N_c \rightarrow \infty$. However, a very large cross-validation set leaves only few data for training thus increasing $G(\hat{\mathbf{w}})$. Obviously, there exists an optimal fraction γ which trades off the conflicting aims. Assume that the quality of the cross-validation estimator is measured by

$$\text{MSE}(\gamma) = E_{\mathcal{D}} \left\{ \left[\hat{G}(\hat{\mathbf{w}}) - G(\hat{\mathbf{w}}) \right]^2 \right\} \quad (5)$$

where $E_{\mathcal{D}}\{\cdot\}$ is the expectation w.r.t. all training data. Further, assume that the loss is the mean square error and that the training data are independent. Since $E_{\mathcal{C}}\{\hat{G}(\hat{\mathbf{w}})\} = G(\hat{\mathbf{w}})$ evaluating Eq. (5) gives

$$\text{MSE}(\gamma) = E_{\mathcal{T}} \left\{ \frac{1}{N_c} \left[E_{\mathcal{C}} \{e^4(\hat{\mathbf{w}})\} - G^2(\hat{\mathbf{w}}) \right] \right\} \quad (6)$$

Using asymptotic expansions (see e.g., [6],[7]) for the terms in Eq. (6) and considering the model to be complete, it is possible to show that the

⁷ $\lceil \cdot \rceil$ denotes rounding upwards to the nearest integer.

optimal fraction is given by $\gamma_{\text{opt}} = 1 - \sqrt{\beta/N}$ where $\beta = 4p\sigma_\varepsilon^2/(1-\xi)$, ξ is the kurtosis of the inherent noise (equal to 3 for Gaussian noise), and $p = \dim(\mathbf{w})$. That is, $\lim_{N \rightarrow \infty} \gamma_{\text{opt}} = 1$ while $N_t = O(\sqrt{N\beta})$ and $N_c = O(N - \sqrt{N\beta})$ asymptotically. It should be emphasized that the choice of γ for a finite small N still needs to be tuned by hand.

The hold-out cross-validation scheme can also be used for comparing generalization errors of different models. Consider the scenario of pruning a nested family of neural net models and suppose that two alternative models with weights $\hat{\mathbf{w}}_1$, $\hat{\mathbf{w}}_2$ both are estimated from \mathcal{T} . If we take $\hat{\mathbf{w}}_2$ to be a subset of $\hat{\mathbf{w}}_1$, i.e., $\dim(\mathbf{w}_2) < \dim(\mathbf{w}_1)$, the hypothesis to be tested is: $G(\hat{\mathbf{w}}_2) > G(\hat{\mathbf{w}}_1)$. Since the models are nested and estimated from the same training set, the corresponding generalization errors are highly dependent. A straight forward procedure which puts error bars on the individual generalization error estimates may fail to unveil the superiority of one model relative to another. The dependence is easily taken into account by analyzing the difference in generalization error, $\Delta\hat{G} = \hat{G}(\hat{\mathbf{w}}_2) - \hat{G}(\hat{\mathbf{w}}_1)$. According to the central limit theorem⁸ $\Delta\hat{G}$ tends to a Gaussian distribution as $N_c \rightarrow \infty$. That is a standard t-test for the hypothesis can be used. If $\Delta\hat{G}/\text{std}(\Delta\hat{G}) < t_\alpha(N_c - 1)$ we reject the hypothesis on an α significance level. $t_\alpha(N_c - 1)$ is the α -fractile of the t -distribution with $N_c - 1$ degrees of freedom, and $\text{std}(\cdot)$ denotes the standard deviation. Define $\Delta e^2(k) = e^2(k, \hat{\mathbf{w}}_2) - e^2(k, \hat{\mathbf{w}}_1)$ then the standard deviation is estimated via $(\widehat{\text{std}}(\Delta\hat{G}))^2 = (N_c - 1)^{-1} N_c^{-1} \sum_{k \in \mathcal{C}} (\Delta e^2(k) - \Delta\hat{G})^2$.

EMPIRICAL GENERALIZATION ERROR DISTRIBUTIONS

We suggest to estimate the generalization error distribution by using leave-out cross-validation [12], [13] and resampling techniques. The basic algorithm is given by:

1. Specify the leave-out fraction γ and determine $N_c = \lceil N\gamma \rceil$. Further specify the number of resamplings $J \leq N!/N_c!(N - N_c)!$.
2. For $j = 1, 2, \dots, J$ split the training set randomly into a cross-validation subset, \mathcal{C}_j , and a training set, $\mathcal{T}_j = \mathcal{D} \setminus \mathcal{C}_j$ not used previously⁹.
3. Train on \mathcal{T}_j with $N_t = N - N_c$ examples to obtain the weight estimate $\hat{\mathbf{w}}_j$ and calculate the empirical mean of the loss on the samples \mathcal{C}_j , which yields the generalization error estimate:

$$\hat{G}_j = \hat{G}_j(\hat{\mathbf{w}}_j) = \frac{1}{N_c} \sum_{k \in \mathcal{C}_j} \ell(y(k), \hat{y}(k); \hat{\mathbf{w}}_j). \quad (7)$$

The training in step 3 can be very time consuming and in [4] we developed an approximate technique for leave-one-out cross-validation.

⁸This also applies when the error signal is a strongly mixing sequence (time-dependent).

⁹Note that this is resampling without replacement, as opposed to the Bootstrap technique.

Ideally, when estimating Eq. (7) we should train and test on independent sets. Moreover, the training sets should be independent. These properties only hold approximately. First, it is very important to stress the significance leaving out a *fraction* γ compared to the standard approach of leaving out a *fixed* number. In the latter case, the different training sets will be too dependent even in the limit of $N \rightarrow \infty$ ¹⁰. However, as discussed in the previous section by letting $\gamma \rightarrow 1$ and if $N_t = O(\nu \log(N))$, $N_c = O(N - \nu \log(N))$, where ν is a constant, all moments of \hat{G}_j converges¹¹. The number of resamplings J should also be allowed to increase towards infinity as N grows. Secondly, for most signal processing problems, time-dependence can especially for small N cause noise in the estimates. However, asymptotically this is no problem since we expect the input signal to be a strongly mixing sequence, i.e., the time-dependence vanishes for large lags.

From the estimates \hat{G}_j in Eq. (7) it is possible to form the *empirical generalization error distribution*

$$P_{\text{emp}}(G) = \frac{1}{J} \sum_{j=1}^J \mu(G - \hat{G}_{(j)}) \quad (8)$$

where $\hat{G}_{(1)} \leq \hat{G}_{(2)} \leq \dots \leq \hat{G}_{(J)}$ is the sample order statistics, and $\mu(G - \hat{G}_{(j)}) = 1$ when $G \geq \hat{G}_{(j)}$, and zero otherwise.

Since $p(G)$ is highly non-Gaussian and long tailed (which is demonstrated experimentally below), the mean and variance are not sufficient for characterizing the shape of $p(G)$. It may consequently be desirable to consider more robust location and dispersion measures which we are able to calculate with $P_{\text{emp}}(G)$ in hand. In general the location of $p(G)$ delivers an estimate of the level of generalization error. The dispersion conveys the fluctuation in generalization error and might suggest if the current number of examples is sufficient for learning the task properly. We consider the following quantities:

Location:

- The average $\text{avr}(G) = \int G p(G) dG$.
- The trimmed average $\text{tavr}(G) = \int_{G_{5\%}}^{G_{95\%}} G p(G) dG$ which reflects the average in which the highest and lowest 5% of the data are excluded.
- The median $\text{med}(G)$ equal to the $\alpha = 50\%$ fractile $G_{50\%}$.

Dispersion:

- The standard deviation $\text{std}(G) = (\int [G - \text{avr}(G)]^2 dG)^{1/2}$.
- The median absolute deviation $\text{mad}(G) = \text{med}(|G - \text{med}(G)|)$.
- The interquartile range $\text{iqr}(G) = G_{75\%} - G_{25\%}$.

¹⁰This is discussed in the literature of the so-called Jackknife estimators, see e.g., [2], [11, Ch. 5.7]

¹¹This is a generalization of what was stated in the previous section for convergence of the second order moment in Eq. (5).

Due to the fact that $p(G)$ follows a χ^2 like distribution, we might consider a transformation of G in order to make it more well behaved. In the general family of Box-Cox transformations (see e.g., [11, Ch. 2.8]) we found that a suitable transformation¹² is $\tilde{G} = \log(1 + G)$.

As in the previous section it is possible to compare the generalization ability e.g, by comparing estimated average generalization errors for two models described by $\mathbf{w}_2, \mathbf{w}_1$. Define the associated estimates $\widehat{avr}(G)_i = J^{-1} \sum_{j=1}^J G_{ij}(\hat{\mathbf{w}}_{ij}), i = 1, 2$, and the difference $\Delta \widehat{avr}(G) = \widehat{avr}(G_2) - \widehat{avr}(G_1)$. For J large $\Delta \widehat{avr}(G)$ tends to a Gaussian distribution by the central limit theorem with standard deviation given by

$$std(\Delta \widehat{avr}(G)) = \sqrt{\frac{1}{(J-1)J} \sum_{j=1}^J [G_{2j} - G_{1j} - \Delta \widehat{avr}(G)]^2} \quad (9)$$

Here the individual differences are assumed to be independent. A standard t-test (as described previously) can then be applied.

NUMERICAL EXAMPLES

Consider the following data generating system: $y(k) = \mathbf{x}^\top(k) \mathbf{w}^\circ + \varepsilon(k)$. $\mathbf{x}(k)$ follows a $L = 10$ variate Gaussian distribution $\mathcal{N}(\mathbf{0}, \mathbf{H})$ with \mathbf{H} chosen as a random positive definite symmetric matrix. $\mathbf{x}(k)$ is time-dependent: each component is a first order AR-process with coefficient 0.6518 scaled to give unit variance; thus implementing a low-pass filter with memory length approx. equal to 7. The noise $\varepsilon(k) \sim \mathcal{N}(0, \sigma_\varepsilon^2)$ is i.i.d. and independent of $\mathbf{x}(k)$. The weights, \mathbf{w}° , were chosen independently from a $\mathcal{N}(0, 1)$ distribution.

We generated $Q = 30000$ independent training sets of size $N = 20$ and trained with a $p = 10$ dimensional linear model using the mean square error cost (without regularization) to obtain the estimates $\hat{\mathbf{w}}^{(i)}, i \in [1; Q]$. This enables a highly accurate estimate of the considered generalization performance measures. As an example, the “true” average generalization error is calculated by $avr(G) = Q^{-1} \sum_{i=1}^Q G(\hat{\mathbf{w}}^{(i)})$ where $G(\hat{\mathbf{w}}^{(i)}) = \sigma_\varepsilon^2 + (\hat{\mathbf{w}}^{(i)} - \mathbf{w}^\circ)^\top \mathbf{H} (\hat{\mathbf{w}}^{(i)} - \mathbf{w}^\circ)$. For $q = 500$ of the $Q = 30000$ training sets we applied the leave-out procedure with $\gamma = 0.25$, $J = 500$ to obtain the weight estimates $\hat{\mathbf{w}}_j^{(i)}$, and corresponding generalization error estimates \hat{G}_j , $j \in [1; J]$. For comparison we calculated the FPE [1] estimate of $avr(G)$ by $FPE^{(i)} = S_N(\hat{\mathbf{w}}^{(i)}) \cdot (N + p)/(N - p)$.

Fig. 1 shows the obtained generalization error probability distributions. Table 1 shows a comparison of the suggested measures of location and dispersion. We consider the transformed variables \tilde{G} which experimentally showed to improve the performance significantly over G . The table indicates that the proposed leave-out technique is fairly accurate for estimating the location and

¹²This implies: $\tilde{G} = 0$ for $G = 0$.

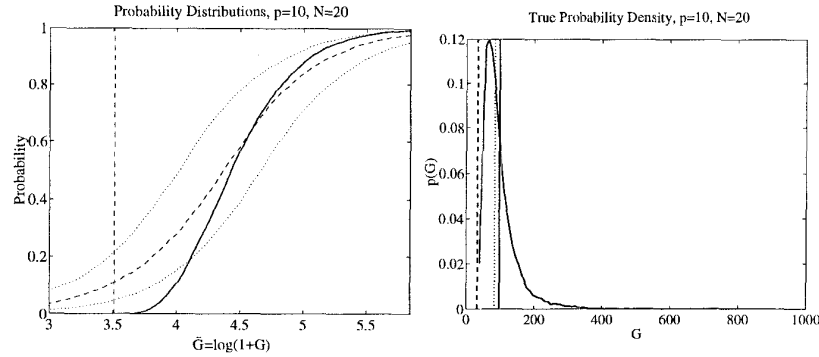


Figure 1: *Left panel:* True (solid) and empirical generalization error distributions as a function of $\tilde{G} = \log(1 + G)$. The dashed line indicates the median distribution of $q = 500$ obtained by leave-out cross-validation while the dotted lines denote the 25% and 75% fractiles. The vertical dashed line is the lower bound $\log(1 + \sigma_\epsilon^2)$. *Right panel:* True generalization error density $p(G)$ estimated from $Q = 30000$ replications. The vertical dashed line is the lower bound $\sigma_\epsilon^2 = 32.45$. The vertical solid line denotes the average, and the vertical dotted the median. Note that $p(G)$ is highly non-Gaussian and long-tailed (ranges to $G = 1000$ approx.). This implies that the classical measure of location, viz. the average overestimates the typical (the mode) generalization error.

dispersion measures even though the number of training data is only twice as large as the number of weights. Definitely, the leave-out method outperforms the classical FPE estimate at the expense of increased computational complexity. However, the framework offers the possibility of estimating other quantities which are not possible in the asymptotic framework on which FPE relies.

We considered furthermore the comparison of two competing linear models: w_1 with dimension $p_1 = 10$, and w_2 with dimension $p_2 = 9$ which consequently is an incomplete model of the true data generating system. The true difference in average generalization ability $\Delta_{avr}(G)$ is positive thus indicating that one should prefer model 1 over model 2. Using the same simulation setup as described above the t-test on a specified $\alpha = 5\%$ significance level resulted in that the hypothesis fails to be accepted in approx. 30% of the cases. More over we considered estimating $\text{Prob}(G_2 > G_1)$ from the empirical distributions. It turned out that the estimate tend to under estimate the probability by 20%. Further, it is somewhat more robust than the estimates of the location measures of the generalization error difference.

¹⁴When considering \tilde{G} the FPE estimate becomes: $\log(1 + \sigma_\epsilon^2) + \sigma_\epsilon^2 p / (1 + \sigma_\epsilon^2) N$ with $\sigma_\epsilon^2 = S_N(\hat{w}) / (N - p)$.

Measure	Min.	25% fract.	Median	75% fract.	Max.
FPE	-62.0	-19.3	-12.3	-5.83	19.7
$\widehat{avr}(G)$	-24.0	-15.5	-7.00	2.27	6.72
$\widehat{tavr}(G)$	-23.5	-15.2	-6.47	2.38	7.05
$\widehat{med}(G)$	-23.2	-15.3	-6.47	3.36	9.31
$\widehat{std}(G)$	44.6	49.9	53.8	62.9	89.4
$\widehat{mad}(G)$	21.6	29.9	40.7	51.4	83.6
$\widehat{iqr}(G)$	20.1	28.6	40.3	53.2	89.4

Table 1: The values are deviations from the true measures in percent when considering the transformed value $\tilde{G} = \log(1 + G)$, e.g., $100\% \cdot (\widehat{avr}(G) - avr(\tilde{G})) / avr(\tilde{G})$. The columns indicate the fluctuation in the deviations w.r.t. the $q = 500$ times the leave-out cross-validation procedure is replicated. As regards FPE¹⁴, the fluctuations are based on $Q = 30000$ replications. In median the location measures $avr(G)$, $tavr(G)$, and $med(G)$ seem to underestimate but are still fairly close to zero, and closer than the estimate of $avr(\tilde{G})$ obtained by FPE. Moreover, the fluctuations are much smaller when considering the fractiles. As regards the dispersion measures $std(G)$, $mad(G)$, and $iqr(G)$ we note that they overestimates by roughly 50%; however, with fairly small amount of fluctuation. The small fluctuation relates strongly to the fact that the transformed variable is used.

CONCLUSION

This paper reports on generalization performance measures which can be attained empirically by using the cross-validation technique in combination with resampling. The major advantage is that the framework provides insight into the shape of the generalization error probability distribution by considering different location and dispersion measures. Traditionally, only the average generalization error has been investigated; however a simple simulation study shows that this measure overestimates the typical generalization performance of a model estimated from a randomly selected set of N examples. Moreover, the assessment of dispersion measures allows for testing the hypothesis whether a model generalizes significantly better than a competitor.

ACKNOWLEDGMENTS

This research was supported by the Danish Natural Science and Technical Research Councils through the Computational Neural Network Center (CONNECT). JL furthermore acknowledge the Radio Parts Foundation for financial support.

REFERENCES

- [1] H. Akaike, "Fitting Autoregressive Models for Prediction," Annals of the Institute of Statistical Mathematics, vol. 21, pp. 243–247, 1969.

- [2] T. Fox, D. Hinkley & K. Larntz, "Jackknifing in Nonlinear Regression," Technometrics, vol. 22, pp. 29–33, 1980.
- [3] S. Geman, E. Bienenstock & R. Doursat, "Neural Networks and the Bias/Variance Dilemma," Neural Computation, vol. 4, pp. 1–58, 1992.
- [4] L.K. Hansen & J. Larsen, "Linear Unlearning for Cross-Validation," submitted for publication, 1995. harlar.luloo.ps.Z is retrieved by anonymous ftp at ei.dtu.dk.
- [5] J. Larsen, "A Generalization Error Estimate for Nonlinear Systems," in S.Y. Kung, F. Fallside, J. Aa. Sørensen & C.A. Kamm (eds.) Neural Networks for Signal Processing 2: Proceedings of the 1992 IEEE-SP Workshop, Piscataway, New Jersey: IEEE, 1992, pp. 29–38.
- [6] J. Larsen, Design of Neural Network Filters, Ph.D. Thesis, Electronics Institute, The Technical University of Denmark, March 1993.
- [7] J. Larsen & L.K. Hansen, "Generalization Performance of Regularized Neural Network Models," in J. Vlontzos, J.-N. Hwang & E. Wilson (eds.), Proceedings of the IEEE Workshop on Neural Networks for Signal Processing IV, Piscataway, New Jersey: IEEE, pp. 42–51, 1994.
- [8] L. Ljung, System Identification: Theory for the User, Englewood Cliffs, New Jersey: Prentice-Hall, 1987.
- [9] J. Moody, "Note on Generalization, Regularization, and Architecture Selection in Nonlinear Learning Systems," in B.H. Juang, S.Y. Kung & C.A. Kamm (eds.) Proceedings of the first IEEE Workshop on Neural Networks for Signal Processing, Piscataway, New Jersey: IEEE, pp. 1–10, 1991.
- [10] N. Murata, S. Yoshizawa and S. Amari, "Network Information Criterion — Determining the Number of Hidden Units for an Artificial Neural Network Model," IEEE Transactions on Neural Networks, vol. 5, no. 6, pp. 865–872, Nov. 1994.
- [11] G.A.F. Seber & C.J. Wild, Nonlinear Regression, New York, New York: John Wiley & Sons, 1989.
- [12] M. Stone, "Cross-validated Choice and Assessment of Statistical Predictors," Journal of the Royal Statistical Society B, vol. 36, no. 2, pp. 111–147, 1974.
- [13] G.T. Toussaint, "Bibliography on Estimation of Misclassification," IEEE Transactions on Information Theory, vol. 20, no. 4, pp. 472–479, July 1974.
- [14] V.N. Vapnik, Estimation of Dependences Based on Empirical Data, New York, New York: Springer-Verlag, 1982.
- [15] V.N. Vapnik, "Principles of Risk Minimization for Learning Theory," in J.E. Moody, S.J. Hanson, R.P. Lippmann (eds.) Advances in Neural Information Processing Systems 4, Proceedings of the 1991 Conference, San Mateo, California: Morgan Kaufmann Publishers, 1992, pp. 831–838.
- [16] H. White, "Learning in Artificial Neural Networks: A Statistical Perspective," Neural Computation, vol. 1, pp. 425–464, 1989.
- [17] H. White, "Consequences and Detection of Misspecified Nonlinear Regression Models," Journal of the American Statistical Association, vol. 76, no. 374, pp. 419–433, June 1981.